

LECTURE 16

WHAT FUNCTIONS $K(x,y)$ ARE KERNELS?

- Can we express \forall function $K(\cdot, \cdot)$ as $K(x,y) = \langle \phi(x), \phi(y) \rangle$?
- No: the matrix $[K(x_i, x_j)]_{i,j=1}^n = [\langle \phi(x_i), \phi(x_j) \rangle]_{i,j=1}^n$ is a Gram matrix \Rightarrow must be symmetric PSD: $K \succeq 0$. (HW 8)
- This necessary condition is also \approx sufficient:

Thm (Mercer) Let $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ be a symmetric, continuous function such that

$$[K(x_i, x_j)]_{i,j=1}^n \succeq 0 \quad \forall x_1, \dots, x_n \in \mathbb{R}^d.$$

Then \exists Hilbert space \mathcal{H} and a map $\phi: \mathbb{R}^d \rightarrow \mathcal{H}$ such that

$$K(x,y) = \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} \quad \forall x,y \in \mathbb{R}^d.$$

↑
"feature space"

↑
"feature map"

⌈ This is an ∞ -dimensional version of the HW8 problem:
A PSD matrix is a Gram matrix ($x \mapsto i, y \mapsto j$) ⌋

- Mercer's condition (PSD) is difficult to check.

However, using it, we can build new kernels from old ones using simple rules (HW), such as addition, exponentiation, etc.

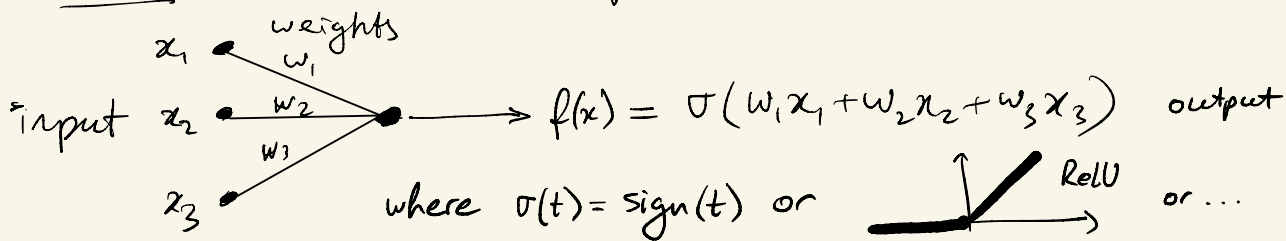
- Examples: (a) RBF $K(x,y) = \exp\left(-\frac{\|x-y\|_2^2}{\sigma^2}\right)$,
(b) polynomial $K(x,y) = (1 + \langle x,y \rangle)^p, \dots$

HW

Kernels arise in:

NEURAL NETWORKS

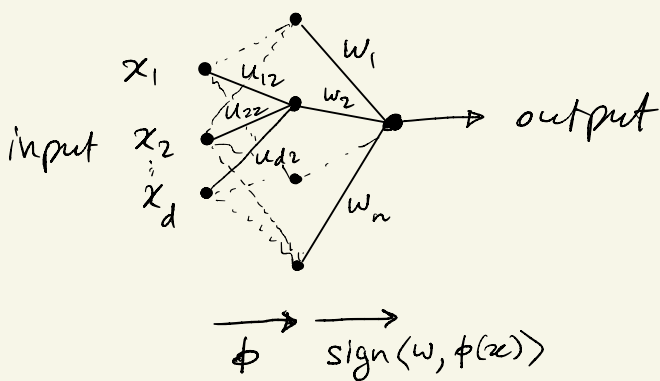
Def A neuron is a composition of linear & nonlinear function



Ex: SVM is a neuron: $\text{sign}(w, x) = \text{sign}(w_1 x_1 + \dots + w_d x_d)$
 We train weights w_j

Def A neural network is a superposition of neurons.

We train weights $(u_{ij}), (w_j)$.



Connection to kernels:

The first layer computes a feature map $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^n$:

$$\phi(x) = \begin{bmatrix} \text{sgn}(u_1, x) \\ \vdots \\ \text{sgn}(u_n, x) \end{bmatrix}$$

Kernel:

$$K(x, y) = \langle \phi(x), \phi(y) \rangle = \frac{1}{n} \sum_{i=1}^n \text{sgn}(u_i, x) \text{sgn}(u_i, y) \xrightarrow{n \rightarrow \infty}$$

Let's initialize all weights $u_{ij} \sim N(0, 1)$ iid

$$\xrightarrow{n \rightarrow \infty} \mathbb{E} \text{sgn}(u, x) \text{sgn}(u, y) \text{ where } u \sim N(0, I_n)$$

$$= \boxed{\frac{2}{\pi} \arcsin \langle x, y \rangle}$$

Hence: a ^{1-layer} neural network \approx kernel SVM.

But: after initialization, a neural network is training weights u_{ij}
 \Rightarrow trains kernel K .

Dynamics of the prediction function is described by the neural tangent kernel (NTK).

A REFRESHER IN LINEAR ALGEBRA

SPECTRAL DECOMPOSITION

SEE MORE IN
KDP-2022

Thm (Spectral thm for symmetric matrices)

\forall symmetric $n \times n$ matrix A ,

(a) all eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_n$ are real;

(b) the unit eigenvectors u_1, \dots, u_n are orthonormal (\Rightarrow o.n. basis)

Cor (Spectral decomposition) \forall symmetric $n \times n$ matrix A :

$$A = \sum_{i=1}^r \lambda_i u_i u_i^T \quad \text{where } r = \text{rank}(A)$$

and (λ_i, u_i) are eigenvalues & unit eigenvectors of A .

SKIP

$$\forall x \in \mathbb{R}^n: x = \sum_{i=1}^n \langle u_i, x \rangle u_i \quad (\text{o.n. basis decomposition})$$

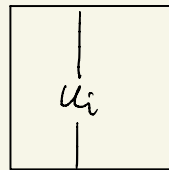
$$\Rightarrow Ax = \sum_{i=1}^n \langle u_i, x \rangle \underbrace{A u_i}_{\lambda_i u_i} = \underbrace{\left(\sum \lambda_i u_i u_i^T \right)}_B x$$

$$\Rightarrow Ax = Bx \quad \forall x \xrightarrow{\text{unr?}} A=B.$$

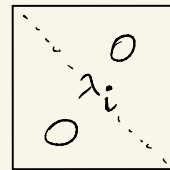
(a) Matrix form:

$$A = U \Lambda U^T$$

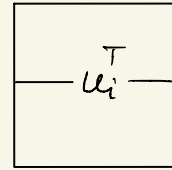
$$A = U \cdot \Lambda \cdot U^T$$



orthog.



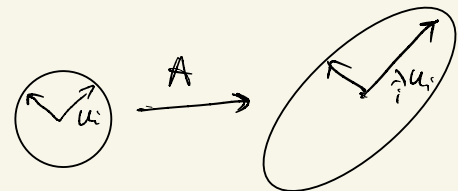
diagonal



orthog.

(b) Geometric form

(c) Optimization form:



$$\lambda_1 = u_1^T A u_1 = \max_{x \text{ unit}} x^T A x, \quad u_1 = \text{argmax} \dots \lambda_1 \text{ times}$$

$$\lambda_2 = u_2^T A u_2 = \max_{\substack{x \perp u_1 \\ \text{unit}}} x^T A x, \quad u_2 = \text{argmax} \dots$$

$$\lambda_3 = u_3^T A u_3 = \max_{x \perp \{u_1, u_2\}} x^T A x, \quad u_3 = \text{argmax} \dots$$

$$\dots$$

$$\lambda_n = u_n^T A u_n = \min_{x \text{ unit}} x^T A x, \quad u_n = \text{argmin} \dots$$